



PATENT ABSTRACTS OF JAPAN

(11)Publication number : 08-305710

(43)Date of publication of application : 22.11.1996

(51)Int.Cl. G06F 17/30

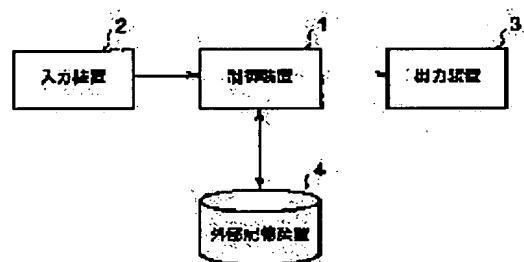
(21)Application number : 07-106582 (71)Applicant : TOSHIBA CORP
TOSHIBA COMPUT ENG CORP
(22)Date of filing : 28.04.1995 (72)Inventor : OZAKI TOSHIHIRO
IWAI ISAMU

(54) METHOD FOR EXTRACTING KEY WORD OF DOCUMENT AND DOCUMENT RETRIEVING DEVICE

(57)Abstract:

PURPOSE: To automatically and effectively extract a keyword to be a key for a document.

CONSTITUTION: This document retrieving device is provided with a control device 1 consisting of a CPU for executing document keyword allocating processing and data processing and a memory, an input device 2 consisting of a keyboard or the like for inputting a document, a processing instruction, etc., an output device 3 consisting of a display or the like for displaying a keyword allocated to a document and processing results, and an external storage device 4 such as an HDD for storing a data base or the like to be used for the allocation of a keyword and constituted so as to compare a certain document with another document stored in a document data base under the control of the control device 1 and extract a characteristic document keyword.



LEGAL STATUS

[Date of request for examination] 16.03.2000

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平8-305710

(43) 公開日 平成8年(1996)11月22日

(51) IntCl.⁶
G 0 6 F 17/30

識別記号

庁内整理番号
9194-5L

F I
G 0 6 F 15/401

技術表示箇所

3 1 0 A

審査請求 未請求 請求項の数10 O L (全 14 頁)

(21) 出願番号 特願平7-106582

(22) 出願日 平成7年(1995)4月28日

(71) 出願人 000003078

株式会社東芝

神奈川県川崎市幸区堀川町72番地

(71) 出願人 000221052

東芝コンピュータエンジニアリング株式会
社

東京都青梅市新町1381番地1

(72) 発明者 尾崎 敏宏

東京都青梅市新町1381番地1 東芝コンピ
ュータエンジニアリング株式会社内

(72) 発明者 岩井 勇

東京都青梅市末広町2丁目9番地 株式会
社東芝青梅工場内

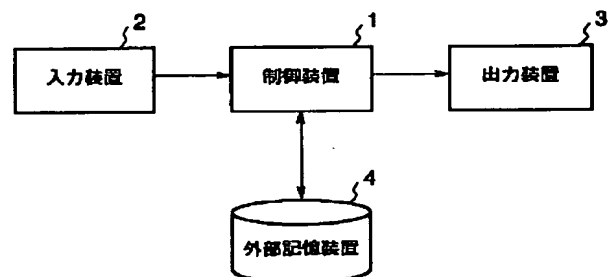
(74) 代理人 弁理士 鈴江 武彦

(54) 【発明の名称】 文書のキーワード抽出方法及び文書検索装置

(57) 【要約】

【目的】 本発明は、文書の鍵となるキーワードを自動的かつ効果的に抽出することを目的としたものである。

【構成】 文書のキーワード付け処理やデータ処理を行なうCPUやメモリからなる制御装置1と、文書や処理指示などを入力するキーボード等からなる入力装置2、文書に付与されたキーワードや処理結果を表示するディスプレイ等からなる出力装置3、文書のキーワード付けのためのデータベースなどを格納するHDD等の外部記憶装置4とを備え、制御装置1の制御の下に、文書データベースの他の文書と比較して特徴のある文書のキーワードを抽出することができる。



【特許請求の範囲】

【請求項 1】 文書に含まれる単語のうち、文書の要旨を知るための鍵となるキーワードを抽出する方法であって、

複数文書から文書に含まれる単語を抽出し、当該の単語を複数文書内での出現回数順に並び替えてランク付けし、

抽出した単語の種数で正規化した総文書単語ランク付けデータベースと、キーワード抽出対象文書に対して前記総文書単語ランク付けデータベースと同様の方法で作成した単語ランク付けデータベースを基にして、キーワード抽出対象文書から抽出した単語に関する総文書単語ランク付けデータベースと単語ランク付けデータベースでの各々のランクを求め、

その差分が許容された範囲以上に単語ランク付けデータベースでのランクが高い場合に、その単語がキーワード抽出対象文書のキーワードであると判断することを特徴とする文書のキーワード抽出方法。

【請求項 2】 総文書単語ランク付けデータベースは、キーワード抽出対象文書の単語ランク付けデータベースの内容を追加して更新が可能なることを特徴とする請求項 1 記載の文書のキーワード抽出方法。

【請求項 3】 複数の文書から抽出された単語毎に、複数文書の各文書に対してその単語の有無を記録した単語存在データベースを保持し、多くの文書に含まれている単語はキーワードとしないことを特徴とする請求項 1 又は 2 記載の文書のキーワード抽出方法。

【請求項 4】 複数の文書から抽出した単語の各々に対して複数の文書内の出現回数をカウントして、その出現回数順に並び替えてランク付けを行ない抽出した単語の種数で正規化した総文書単語ランク付けデータベースを作成し、このデータベースを用いて文書のキーワードの抽出を行なう文書のキーワード付け装置であって、キーワード抽出対象文書内の文字列を単語に分割する文書—単語分割手段と、文書から分割した単語が文書中に何回出現するかをカウントする単語出現回数カウント手段と、

文書内の単語の出現回数で単語を並び替えてランク付けし抽出した単語の種数で正規化した単語ランク付け手段と、

文書から分割した単語について、総文書単語ランク付けデータベースと単語ランク付け手段とから夫々に得られたランクを比較して、ランクの差が指定された許容範囲を越える程度にキーワード抽出対象文書内に出現する頻度が高い場合に当該文書のキーワードと判断するキーワード判別手段と、

キーワードやその他の処理結果を出す出力手段とを具備したことを特徴とする文書のキーワード付け装置。

【請求項 5】 総文書単語ランク付けデータベースは、キーワード抽出対象文書の単語ランク付けデータベース

の内容を追加して更新が可能なることを特徴とする請求項 4 記載の文書のキーワード付け装置。

【請求項 6】 キーワード判別手段は、着目単語がキーワード抽出対象のキーワードか否かの判断をする際に、複数の文書から抽出された単語毎に、複数文書の各文書に対してその単語の有無を記録した単語存在データベースを保持し、多くの文書に含まれている単語はキーワードとしない単語存在データベースを用いて、どの文書にも現れるような単語をキーワードとしないことを特徴とする請求項 4 記載の文書のキーワード付け装置。

【請求項 7】 文書を分野に分類する文書分類装置で、分類対象の分野毎にその分野を代表するキーワードを付与しておいた分野別キーワード情報を用いて文書を分野に分類する文書分類装置であって、

文書のキーワードを抽出するキーワード抽出手段と、文書のキーワードと分野別キーワードを比較して一致した分野を当該文書の所属する分野であると判断する分野判断手段と、

文書を入力する入力手段と、結果を出力する出力手段とを具備したことを特徴とする文書分類装置。

【請求項 8】 キーワード抽出手段は、請求項 1 記載の文書のキーワード抽出方法方法であることを特徴とする請求項 7 記載の文書分類装置。

【請求項 9】 検索キーワードを含む文書をフルテキストサーチによって検索する文書検索装置であって、検索キーの入力や処理指示を行なう入力手段と、検索手段と、

検索結果を出力する出力手段と、

検索手段によって得られた検索キーワードを含む文書について検索キーワードの文書内での重要度によって並び替えを行って最終的な検索結果とする文書重要度判別手段とを具備することを特徴とする文書検索装置。

【請求項 10】 文書重要度判別手段は、検索の実行前に予め作成しておいた単語の種数で正規化した総文書単語ランク付けデータベース、及び単語ランク付けデータベースを基にして、2つのデータベースから得られる検索キーワードのランクの差を検索手段で求められた文書毎に数値化して重要度を判別することを特徴とする請求項 9 記載の文書検索装置。

40 【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、文書の要旨を知るための重要な単語をキーワードとして抽出する作業を支援し、文書を分類する作業を支援すると共に、文書データベース等に登録された文書をキーワードをもとに検索するようにした、文書のキーワード抽出方法と、キーワード付け装置、文書分類装置、及び文書検索装置に関する。

【0002】

50 【従来技術】 従来、ある文書が何について述べている

かを分かりやすくするための手法には様々なものがある。一般的には、本文の要約や表題などで、また新聞等では、見出しを付与して内容を簡潔かつ効果的に伝える手法がとられている。また、文書の要旨を知るための鍵となる単語（キーワード）を取り出して、文書に付与する手法も存在する。

【0003】文書にキーワードを付与する方法は、文書データベースに採用されている。文書の文書データベースへの登録時に、付与したキーワードも同時に登録する。この付与したキーワードは、文書データベースの検索時の検索キーワードとして用いられる。従来、このキーワードの付与作業は人手によって行なわれてきた。

【0004】一方、機械的に文書中の単語を抽出するために、日本語解析や形態素解析を用いる方法も存在するが、この方法のみでは、抽出した単語が文書の要旨を知る鍵となるキーワードであるかの識別をつけることはできず、これを文書のキーワードとするには難点があった。

【0005】また、文書が何の分野について述べられているかを判断し分類する作業については、効果的な分類を実施するためには文書の内容を理解しなければならず、これも人手によって分類が行なわれている場合が多い。

【0006】また、文書の全文を検索対象としたフルテキストサーチでは、文書データベースに対して検索キーワードを含む文書を検索して得られた結果には、検索キーワードに関連が深い文書だけでなく、その検索キーワードが引用されているだけの相対的に重要度も低い文書も含まれていた。

【0007】

【発明が解決しようとする課題】文書が何について書かれているものであるかを、文書の要旨を知るための鍵となるキーワードの形にする作業は人手によるものが多く、また、内容を理解するために、かなりの時間と専門性を必要とする。

【0008】また、従来の機械的な文書からの単語の切り出しのみでは、文書のキーワードを抽出することができなかった。本発明は、上記事情を考慮して成されたものであり、文書の鍵となるキーワードを自動的かつ効果的に抽出することを目的としたものである。

【0009】また、本発明は、これまで人手によって分類していた文書を前記基本発明である文書のキーワードの自動付与技術を用いて分類し、文書の分類作業の軽減を目的とする。

【0010】また、フルテキストサーチでの検索支援があるが、文書の検索キーワードに対する有無を検索結果として出力するだけでなく、検索結果を検索キーの重要度が高いものについて並び替えて出力を行ない、ユーザーの要求する文書が探しやすくなる様な検索作業の効率化を目的とするものである。

【0011】

【課題を解決するための手段】本発明は上記目的を達成するため、文書に含まれる単語の内、文書の要旨を知るための鍵となるキーワードを抽出する方法であって、複数文書から文書に含まれる単語を抽出し、当該の単語を複数文書内での出現回数順に並び替えてランク付けし、抽出した単語の種数で正規化した総文書単語ランク付けデータベースと、キーワード抽出対象文書に対して前記総文書単語ランク付けデータベースと同様な方法で作成した単語ランク付けデータベースを基にして、キーワード抽出対象文書から抽出した単語に関する総文書単語ランク付けデータベースと単語ランク付けデータベースでの各々のランクを求め、その差分が許容された範囲以上に単語ランク付けデータベースでのランクが高い場合に、その単語がキーワード抽出対象文書のキーワードであると判断することを特徴とする文書のキーワード抽出方法にある。

【0012】また、本発明は上記目的を達成するため、複数の文書から抽出した単語の各々に対して複数の文書内の出現回数をカウントして、その出現回数順に並び替えてランク付けを行ない抽出した単語の種数で正規化した総文書単語ランク付けデータベースを作成し、このデータベースを用いて文書のキーワードの抽出を行なう文書のキーワード付け装置であって、キーワード抽出対象文書内の文字列を単語に分割する文書—単語分割手段と、文書から分割した単語が文書中に何回出現するかをカウントする単語出現回数カウント手段と、文書内の単語の出現回数で単語を並び替えてランク付けし抽出した単語の種数で正規化した単語ランク付け手段と、文書から分割した単語について総文書単語ランク付けデータベースと単語ランク付け手段とから夫々に得られたランクを比較してランクの差が指定された許容範囲を越える程度にキーワード抽出対象文書内に出現する頻度が高い場合に当該文書のキーワードと判断するキーワード判別手段と、キーワードやその他の処理結果を出す出力手段とを具備したことを特徴とする文書のキーワード付け装置にある。

【0013】また、本発明は上記目的を達成するため、文書を分野に分類する文書分類装置で、分類対象の分野毎にその分野を代表するキーワードを付与しておいた分野別キーワード情報を用いて文書を分野に分類する文書分類装置であって、文書のキーワードを抽出するキーワード抽出手段と、文書のキーワードと分野別キーワードを比較して一致した分野を当該文書の所属する分野であると判断する分野判断手段と、文書を入力する入力手段と、結果を出力する出力手段とを具備したことを特徴とする文書分類装置にある。

【0014】更に、本発明は上記目的を達成するため、検索キーワードを含む文書をフルテキストサーチによって検索する文書検索装置であって、検索キーの入力や処

理指示を行なう入力手段と、検索手段と、検索結果を出力する出力手段と、検索手段によって得られた検索キーワードを含む文書について検索キーワードの文書内での重要度によって並び替えを行って最終的な検索結果とする文書重要度判別手段とを具備することを特徴とする文書検索装置にある。

【0015】

【作用】上記構成に於いては、複数の文書から形態素解析等で単語切りを行って抽出した単語の複数の文書に対する出現回数の順にランク付けを行ない出現した単語の種類の総数で正規化を行った基準となる総文書単語ランク付けベースと、キーワード抽出対象文書に対して前記の場合と同様に作成した単語ランク付けデータベースとを用いて、着目単語の2つのデータベースのランクの差を利用して対象文書のキーワードか否かの判定を行なうことにより、単純に文書を単語切りして抽出した単語がキーワードとはならず、また、2つのデータベースを相対的に参照するために、文書内には頻繁に出現するが重要度が低いと判断される単語、例えば特許明細書における「手段」という単語は、対象文書のキーワードとして判定されることは無くなり、キーワード抽出対象文書の要旨を知る鍵となる的確なキーワードを抽出することが可能となる。

【0016】また、上記構成に於いては、予め分野を代表するキーワードを分類したい分野毎に分野別キーワードとして付与しておくことで、分類対象文書から抽出したキーワードと一致する分野別キーワードを比較することが可能となり、一致した分野別キーワードを持つ分野に、分類対象文書を分類することが機械的に可能となり、人手を省略した効率的な分類が実現できる。

【0017】更に、上記構成に於いては、予め作成しておいた総文書単語ランク付けデータベースと検索対象の各文書毎に作成した単語ランク付けデータベースを準備しておくことで、検索キーワードを含む文書の検索キーワードの重要度が数値的に判断可能となり、検索結果として検索キーワードの重要度に並べ替えを行った検索結果がユーザーに提示可能となり、検索の効率を向上することができる。

【0018】

【実施例】本発明の概要は次の通りである。

(1) 文書に含まれる単語が重要であるかの識別を行なう情報として、文書の鍵となる重要な単語は文書中に多く現れるという前提に着目した。また、文書中に多く現れるものの、さして重要ではない単語をキーワードとする事を防ぐため、文書データベース全体の情報を基にして、相対的に文書のキーワードを抽出する。詳しく述べるならば、複数の文書に現れる単語をその出現回数によってランク付けして単語種の総数で正規化した情報と、キーワードを抽出しようとする文書に現れる単語をその出現回数によってランク付けして単語種の総数で正規化

した情報の2つの情報からそれぞれにランクを取り出し、キーワードであるか否かを判断する際にそのランクの差分が誤差範囲と見なせない程、単語が数多く出現すると判断された場合に、当該文書のキーワードとする手法を採る。

【0019】また、本発明は、文書を入力したり処理指示を与える入力手段と、文書を単語に分割する文書単語分割手段と、文書から分割した単語が文書中に何回出現するかをカウントする単語出現回数カウント手段と、文書内の単語の出現回数を用いて計算した結果を基に単語をランク付けして単語種の総数で正規化した単語ランク付け手段と、複数の文書に対して文書中に含まれる単語毎に出現回数の総計を記録して出現回数を用いて計算した結果を基に単語をランク付けした総文書単語ランク付けデータベースと、分割した文書中の単語がキーワードとなる単語かを判断するキーワード判別手段と、キーワードやその他の処理結果を出す出力手段とを具備したことを特徴とする。

【0020】また、前記総文書単語ランク付けデータベースは、キーワード抽出対象文書の単語ランク付けデータベースの内容を追加して更新が可能なことを特徴とする。また、前記キーワード判別手段は、検索実行前に予め作成しておいた、複数の文書から抽出された単語毎に各文書に対してその単語の有無を記録した単語存在データベースを保持し、多くの文書に含まれている単語はキーワードとしないことを特徴とする。

(2) 文書を分野に分類する発明は、分類対象の分野毎にその分野を代表するキーワードを付与しておいた分野別キーワード情報と、文書のキーワードを抽出するキーワード抽出手段と、文書のキーワードと分野別キーワードを比較して一致した分野を当該文書の所属する分野であると判断する分野判断手段と、文書を入力する入力手段と、結果を出力する出力手段とを具備したことを特徴とする。

【0021】また、前記のキーワード抽出手段は、

(1) の方法であることを特徴とする。

(3) 文書検索を支援する発明は、検索キーの入力や処理指示を行なう入力手段と、検索手段と、検索結果を出力する出力手段と、検索手段によって得られた検索キーワードを含む文書について検索キーワードの文書内での重要度によって並び替えを行なって最終的な検索結果とする文書重要度判別手段とを具備したことを特徴とする。

【0022】また、前記の文書重要度判別手段は、検索の実行前に予め作成しておいた(1)記載の総文書単語ランク付けデータベース、及び単語ランク付けデータベースを基にして、2つのデータベースから得られる検索キーワードのランクの差を検索手段で求められた文書毎に数値化して重要度を判別することを特徴とする。

【0023】以下図面を参照して本発明の一実施例を説

明する。

(第1実施例) 本実施例は、文書のキーワード付け装置に関わるものである。

【0024】図1は、文書のキーワード付け装置に係わる概略構成を示すブロック図であり、1は文書のキーワード付け処理やデータ処理を行なうCPU、メモリ等からなる制御装置である。

【0025】2は文書や処理指示などを入力するキーボード等からなる入力装置である。3は文書に付与されたキーワードや処理結果を表示するディスプレイ等からなる出力装置である。

【0026】4は文書のキーワード付けのためのデータベースなどを格納するHDD等の外部記憶装置である。図2は、上記図1に示した制御装置1の詳細例を示したブロック図である。

【0027】制御装置1は初期化部11、入力部12、出力部13、制御部14、文書-単語分割部15、単語出現回数カウンタ部16、単語ランク付け部17、単語ランク積算部18、キーワード判別部19等の制御系と、文書格納バッファ20、単語分割文書格納バッファ21、単語登録バッファ22、出現回数格納バッファ23、キーワード格納バッファ24等の記憶系と、メモリ上に確保した単語ランク付けデータベース25、外部記憶装置に格納した総文書単語ランク付けデータベース26等のデータベース系とから構成されている。

【0028】初期化部11は、記憶系の各バッファの初期化を行なう。入力部12は入力装置2からデータベース作成モードか、キーワード抽出モードかのどちらかのモードの指定と文書の内容等の情報を入力する。出力部13は入力部で指定されたモードがキーワード抽出モードだった時に、その入力された文書のキーワード等の情報を表示装置3に出力する。制御部14は、制御系全体を制御して、データベース作成やキーワードの抽出などの各処理を総合的に制御する。

【0029】入力部12を通して入力されたキーワード抽出対象文書（以下文書と称する）は制御部14を介して文書格納バッファ20に格納される。文書-単語分割部15では、日本語解析などの技術を用いて文書格納バッファ20に格納された文書を単語に分割し、単語と単語の境がわかるようにして単語分割文書格納バッファ21に格納し、さらに文書内で分割された単語を、重複することなく単語登録バッファ22に格納する。

【0030】単語出現回数カウンタ部16は、単語登録バッファ22に登録された単語毎に、該単語が単語分割文書格納バッファ21における出現回数をカウントし、その結果を出現回数格納バッファ23に格納する。単語ランク付け部17はキーワード抽出モードの場合に動作し、出現回数格納バッファ23の単語をその出現回数順に並び替え、出現回数を文書に出現した単語種の総数で正規化したものを単語ランク付けデータベース25に格

納する。

【0031】単語ランク積算部18は、データベース作成モードの場合に動作し、出現回数格納バッファ23に出現した単語の内、未登録の単語を総文書単語ランク付けデータベース26に格納し、登録済みの単語の場合は出現回数を加算して出現回数のデータを更新した後、該データベースに登録されている単語種の総数で再度正規化を行なって外部記憶装置4に格納する。

【0032】キーワード判別部19はキーワード抽出モードの場合に動作し、単語ランク付けデータベース25に格納された単語の各々に対して該データベースと総文書単語ランク付けデータベース26のランクを比較し、その差が閾値を越える程度に出現回数が多いと判断された場合に、該単語をキーワード格納バッファ24に格納する。キーワード格納バッファ24に格納されたキーワードは制御部14を介して出力装置3に出力される。

【0033】次に、本文書のキーワード付け装置の動作と処理の流れについて、図3に示すフローチャートを参照しながら説明する。まず、初期化部11がステップ301にて起動し、各バッファを初期化する。次に、ステップ302にて、入力部12によってモード（キーワード抽出モード、又はデータベース作成モード）の指示を受け、さらに入力部12を介して入力された文書を、制御部14が文書格納バッファ20に格納する。図4は、文書格納バッファ20に格納した文書の例を示した図である。

【0034】文書格納バッファ20に格納された文書は、ステップ303にて起動した文書-単語分割部15によって単語単位に区切られ、区切り符号と共に単語分割文書格納バッファ21に格納される。

【0035】また、ステップ304では文書-単語分割部15によって、抽出した単語を単語登録バッファ22に格納する。図5は、単語分割文書格納バッファ21の例であり、単語と単語の間にはスラッシュを区切り符号として用いてある。図6は、単語登録バッファ22の例であり、文書から抽出した単語の例を示した図である。

【0036】ステップ305では単語出現回数カウンタ部16が起動して、単語分割文書格納バッファ21と単語登録バッファ22から、その抽出した単語毎に出現回数をカウントして、出現回数格納バッファ23に格納する。図7は出現回数格納バッファ23の例であり、単語の文字列とその出現回数を対応づけて格納してある。

【0037】ステップ306ではステップ302で指定されたモードによって処理が異なる。キーワード抽出モードである場合はステップ307へ、データベース作成モードの場合はステップ312へ制御を移す。

【0038】ステップ307では、単語ランク付け部17が起動し、出現回数バッファ23のデータを基にして、各単語を出現回数順に並び替え、さらに文書に出現した単語種の総数で正規化した値を付与して、単語ラン

ク付けデータベース25として格納する。図8は、単語ランク付けデータベース25の例であり、単語の文字列と出現回数と正規化値を対応づけて出現回数順に格納したものを示した図である。

【0039】ステップ308～311ではキーワード判別部19によって処理が行われる。ステップ308では、単語ランク付けデータベースの単語について、単語ランク付けデータベース25と総文書単語ランク付けデータベースの該単語の正規化値を取り出し、この正規化値の差を求める。ステップ309では、ステップ308で求めた差が誤差の範囲を越える程に出現頻度が高いか否かを判断し、キーワードと判断された場合にはステップ310へ、そうでない場合はステップ311に制御を移す。

【0040】ステップ310では、引き続きキーワード判別部19によって処理が行われ、ステップ309でキーワードと判断された単語をキーワード格納バッファ24に格納する。格納したキーワードは制御部14を介して出力部13に送られ、表示が行われる。

【0041】ステップ311では、文書中の単語の全てについてキーワードの判定を行なったか否かを判断し、全てについて判定を行っていない場合はステップ308に制御を移し、そうでない場合は処理を終了する。図9は、キーワード格納バッファ24の例であり、抽出したキーワードを順に格納した図である。

【0042】ステップ312～313はキーワード抽出モードの場合に、単語ランク付け積算部18によって処理が行われる。ステップ312では出現回数格納バッファ23に格納された単語の内、総文書単語ランク付けデータベース26にない単語は、単語ランク付けデータベース25の形で、登録済みの単語は単語ランク付けデータベース25の出現回数を総文書単語ランク付けデータベース26に加算する。ステップ313では、総文書単語ランク付けデータベース26に出現する全ての単語を出現回数で並び替えを行ない、総単語種数で正規化を実施し、その正規化値とともに総文書単語ランク付けデータベース26に格納して処理を終了する。図10は、総文書単語ランク付けデータベース26の例であり、単語ランク付けデータベース25と同じフォーマットで格納したデータを示している。

(第2実施例) 本実施例は文書の分類装置に関わるものである。

【0043】文書分類装置に関わる概略構成を示すブロック図は第1実施例の図1と同様であり、1は文書の分類やデータ処理を行なうCPU、メモリ等からなる制御装置、2は文書や処理指示などを入力するキーボード等からなる入力装置、3は文書分類の結果を表示するディスプレイ等からなる出力装置、4は文書分類のための分野毎のキーワード等を格納するHDD等の外部記憶装置である。

【0044】図11は、図1に示した制御装置1の詳細例を示したブロック図である。制御装置1は、初期化部111、入力部112、出力部113、制御部114、キーワード抽出部115、文書分野判断部116等の制御系と、文書格納バッファ117、抽出キーワード格納バッファ118、分野別キーワード格納バッファ119等の記憶系と、外部記憶装置に格納した分野別付帯キーワード情報120等のデータベース系とから構成されている。

10 【0045】初期化部111は、記憶系の各バッファの初期化を行ない、外部記憶装置4に格納された分野別付帯キーワード情報120を分野別キーワード格納バッファ119に格納する。

【0046】入力部112は入力装置2から入力された文書を制御部114に渡す。出力部113は制御部114から渡された文書の分野等の処理結果を出力する。制御部114は文書分類装置の制御系全体を制御したり、入力部112からの文書データをバッファに格納する等の各処理を総合的に制御する。

20 【0047】キーワード抽出部115は、文書格納バッファ117に格納された文書から文書のキーワードを抽出し、抽出キーワード格納バッファ118に格納する。文書分野判断部116は、分野別キーワード格納バッファ119に格納された分野毎に予め登録されたキーワードと、抽出キーワード格納バッファ118に格納されたキーワードを比較し、一致したキーワードがある場合には該当分野を制御部114に渡す。

【0048】次に、本文書の文書分類装置の動作と処理の流れについて、図12に示すフローチャートを参照しながら説明する。ステップ1201では、初期化部111が動作を行ない、各バッファの初期化、及び外部記憶装置4の分野別付帯キーワード情報120を分野別キーワード格納バッファ119に展開する。図13は、分野別付帯キーワード情報120と分野別キーワード格納バッファ119の格納例を示したもので、分野の名称とその分野に付帯するキーワードを列挙して格納した図である。

【0049】次にステップ1202では、入力部が入力装置2から入力された文書を制御部114に渡し、制御部114は文書格納バッファ117に文書を格納する。文書格納バッファ117は第1実施例で示したものと同じで、この例は図4に示したものと同じである。

【0050】ステップ1203ではキーワード抽出部115が動作を行ない、文書格納バッファ117に格納された文書のキーワードを抽出し、抽出キーワード格納バッファ118に格納する。抽出キーワード格納バッファ118は第1実施例で示したものと同じで、この例は図6で示したものと同じである。

50 【0051】ステップ1204～1207までは文書分野判断部116が動作を行なう。ステップ1204で

は、分野別キーワード格納バッファ119に格納した予め登録してあるキーワードと抽出キーワード格納バッファ118に格納したキーワードの調査を行ない、ステップ1205で一致しているものと判断した場合にはステップ1206に制御を移し、そうでない場合にはステップ1207に制御を移す。

【0052】ステップ1206では一致したキーワードのある分野を制御部114に渡し、制御部114ではキーワードを出力部113に渡して出力装置3に対して出力を行なう。

【0053】ステップ1207では、全ての分野について抽出したキーワードと一致するか否かのチェックを行ない、全ての分野についてチェックを終了していない場合にはステップ1204に制御を移し、そうでない場合には処理を終了する。

(第3実施例) 本実施例は文書の検索装置に関わるものである。

【0054】文書検索装置に係わる概略構成を示すブロック図は第1実施例の図1と同様であり、1は文書の検索や重要度を判断する制御を行なうCPU、メモリ等からなる制御装置、2は処理指示などをを入力するキーボード等からなる入力装置、3は文書や検索の結果を表示するディスプレイ等からなる出力装置、4は検索対象文書や重要度を判別するデータ等を格納するHDD等の外部記憶装置である。

【0055】図11は、図1に示した制御装置1の詳細例を示したブロック図である。制御装置1は、初期化部1401、入力部1402、出力部1403、制御部1404、検索部1405、重要度判別部1406等の制御系と、検索中間結果格納バッファ1407、検索結果格納バッファ1408等の記憶系と、外部記憶装置に格納した検索対象文書1409、重要度判別データ等のデータベース系とから構成されている。

【0056】初期化部1401は、記憶系の各バッファの初期化を行なう。入力部1402は入力装置2から入力された検索の指示や検索キーワード等を制御部1404に渡す。出力部1403は制御部1404から渡された検索結果等の処理結果を出力する。制御部1404は文書検索装置の制御系全体を制御したり、検索結果を出力装置に渡したり等の各処理を総合的に制御する。

【0057】検索部1405は、外部記憶装置4に格納された検索対象文書1409から入力された検索キーワードを含む文書を検索し、該当文書に対応づけられた文書番号を検索中間結果格納バッファ1407に格納する。

【0058】重要度判別部1406では、検索中間結果格納バッファ1407に格納された文書番号の文書について、外部記憶装置4に格納された重要度判別データ1410を用いて重要度のランク付けを行ない、重要度の順番に検索結果格納バッファ1408に文書番号を格納

する。

【0059】次に、本文書の文書検索装置の動作と処理の流れについて、図15に示すフローチャートを参照しながら説明する。ステップ1501では、初期化部1401が動作を行ない、各バッファの初期化を行なう。

【0060】次にステップ1502では入力部1402が動作を行ない、入力装置2から入力された検索キーワードを制御部1404に渡す。ステップ1503では検索部1405が動作を行ない、制御部1404から渡された検索キーワードを含む文書を、検索対象文書1409から検索を行ない、該当文書の文書番号を検索中間結果格納バッファ1407に格納する。図16は、検索中間結果格納バッファ1407の例であり、検索キーワードを含む文書の文書番号が格納されていることを示す図である。

【0061】ステップ1504では重要度判別部1406が動作を行ない、検索中間結果格納バッファ1407に格納された文書番号の文書の重要度を重要度判別データ1410を用いて判断を行ない、文書の重要度の順に文書番号を検索結果格納バッファ1408に格納する。図17は検索結果格納バッファ1408の例であり、検索キーワードを含む文書が重要度の順番に格納されていることを示す図である。

【0062】ステップ1505では制御部1404が動作を行ない、検索結果格納バッファ1408に格納された文書番号を出力部1403に渡し、出力部1403は出力装置3へ文書番号の出力を行なう。

【0063】ステップ1506では制御部1404が動作を行ない、入力部1402から入力された検索続行の指示の有無を基に制御を移行する。検索を続行する場合にはステップ1502に制御を移行し、そうでない場合には処理を終了する。

(第4実施例) 本実施例は、文書のキーワード付け装置に関わるものである。

【0064】本実施例のブロック図、制御装置の詳細ブロック図は第1実施例のものと同じであり、それぞれ、図1と図2が対応する。本実施例と第1実施例との差異は、図3を基にして述べるならば、ステップ306のキーワード抽出モードか否かの判断の後、ステップ307のキーワード抽出対象文書に現れる単語のランク付けのブロックの前に、ステップ312、ステップ313で述べた処理が挿入されるのみである。

(第5実施例) 本実施例は、文書のキーワード付け装置に関わるものである。

【0065】本実施例のブロック図は第1実施例のものと同じであり、図1が対応する。本実施例と第1実施例との差異は、図2を基にして述べるならば、総文書単語ランク付けデータベース26を格納した外部記憶装置4に、複数の文書から抽出された単語毎に、複数文書の各文書に対してその単語の有無を記録した単語存在データ

ベースを保持することである。図18は、単語存在データベースの格納例であり、行方向を単語文字列、列方向を文書番号として、マトリクス上に単語の文書に対する存在の有無を1、0で表した図である。

【0066】この際、キーワード判別部19は、実施例1で述べた、単語ランク付けデータベース25に格納された単語の各々に対して該データベースと総文書単語ランク付けデータベース26のランクを比較し、その差が閾値を越える程度に出現回数が多いと判断する処理に加えて、単語存在データベースを参照し多くの文書に含まれている単語はキーワードとしない処理を行ない、その結果をキーワード格納バッファ24に格納する。

【0067】また、処理の流れについて、第1実施例との差異を図3を基にして述べるならば、ステップ309のキーワードか否かの判断の後、ステップ310のキーワードの出力の前に、単語存在データベースを参照し多くの文書に含まれている単語はキーワードとしない処理が挿入される。

(第6実施例) 本実施例は、文書分類装置に関わるものである。

【0068】本実施例のブロック図は第2実施例のものと同じであり、図1が対応する。本実施例と第2実施例との差異は、図11を基にして述べるならば、キーワード抽出部115が、第1実施例、第4実施例、第5実施例で述べた文書のキーワード抽出方法にて動作することである。

【0069】また、図12での第2実施例との差異は、ステップ1203での文書のキーワードを抽出する方法が、第1実施例、第4実施例、第5実施例で述べた文書のキーワード抽出方法にて処理が行われることである。

(第7実施例) 本実施例は、文書検索装置に関わるものである。

【0070】本実施例のブロック図は第3実施例のものと同じであり、図1が対応する。本実施例と第3実施例との差異は、図14を基にして述べるならば、重要度判別部1406が、第1実施例、第4実施例、第5実施例で述べた文書のキーワード抽出方法にて動作することである。

【0071】また、図15での第3実施例との差異は、ステップ1504での重要度を判別する方法が、第1実施例、第4実施例、第5実施例で述べた文書のキーワード抽出方法にて処理が行われることである。

【0072】

【発明の効果】以上詳記したように本発明によれば、文書データベースの他の文書と比較して特徴のある文書のキーワードを抽出することができる。また、本発明によれば、文書データベースの総文書単語ランク付けデータベースにキーワード抽出対象文書のデータを追加することで、登録時に総文書単語ランク付けデータベースの更新を行なうことができ、またデータ量が増えることによ

り、更に効果的なキーワードの抽出を行なうことができる。

【0073】また、本発明によれば、単語存在データベースによってデータベース内の他の文書に多く現れる単語をキーワードとすることを防ぐことができる。また、本発明によれば、文書の分野分類作業を人手によって行なう必要がなくなり、効率の改善になる。

【0074】また、本発明によれば、文書の分類作業を機械的に行なうため、文書の分類作業の性能が向上する。更に、本発明によれば、検索キーワードを含む文書の検索結果をユーザーに提示する際に、機械的に重要な文書をランクづけるために、ユーザーの検索の効率を向上させることができる。

【図面の簡単な説明】

【図1】本発明の実施例に係る装置の概略構成を示すブロック図。

【図2】本発明の第1実施例に於ける、文書のキーワード付け装置での図1の制御装置1の詳細な構成を示すブロック図。

20 【図3】本発明の第1実施例に於けるシステムの動作の概要を示すフローチャート。

【図4】本発明の第1実施例に於ける文書格納バッファ20と第2実施例の文書格納バッファ117の格納例を示す図。

【図5】本発明の第1実施例に於ける単語分割文書格納バッファ21の格納例を示す図。

【図6】本発明の第1実施例に於ける単語登録バッファ22の格納例を示す図。

30 【図7】本発明の第1実施例に於ける出現回数格納バッファ23の格納例を示す図。

【図8】本発明の第1実施例に於ける単語ランク付けデータベース25の格納例を示す図。

【図9】本発明の第1実施例に於けるキーワード格納バッファ24と第2実施例の抽出キーワード格納バッファ118の格納例を示す図。

【図10】本発明の第1実施例に於ける総文書単語ランク付けデータベース26の格納例を示す図。

【図11】本発明の第2実施例に於ける文書分類装置での図1の制御装置1の詳細な構成を示すブロック図。

40 【図12】本発明の第2実施例に於けるシステムの動作の概要を示すフローチャート。

【図13】本発明の第2実施例に於ける分野別付帯キーワード情報120と分野別キーワード格納バッファ119の格納例を示す図。

【図14】本発明の第3実施例に於ける文書検索装置での図1の制御装置1の詳細な構成を示すブロック図。

【図15】本発明の第3実施例に於けるシステムの動作の概要を示すフローチャート。

50 【図16】本発明の第3実施例に於ける検索中間結果格納バッファ1407の格納例を示す図。

15

【図 17】本発明の第 3 実施例に於ける検索結果格納バッファ 1408 の格納例を示す図。

【図 18】本発明の第 5 実施例に於ける単語存在データベースの格納例を示す図。

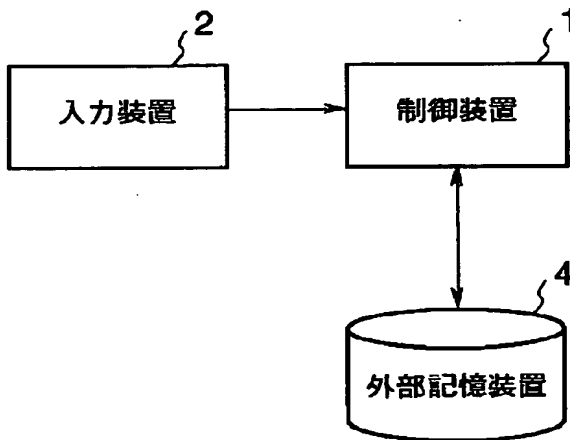
【符号の説明】

1…制御装置、2…入力装置、3…出力装置、4…外部記憶装置、11…初期化部、12…入力部、13…出力部、14…制御部、15…文書-単語分割部、16…

16

単語出現回数カウント部、17…単語ランク付け部、18…単語ランク積算部、19…キーワード判別部、111…初期化部、112…入力部、113…出力部、114…制御部、115…キーワード抽出部、116…文書分野判断部、1401…初期化部、1402…入力部、1403…出力部、1404…制御部、1405…検索部、1406…重要度判別部。

【図 1】



【図 6】

【図 16】

この発明は、半導体に…	文書番号
	7
	24
	39
	146
	255

【図 4】

【図 5】

【図 7】

この発明は、半導体に関するもので、AのB手段を具備する。…

この発明は、半導体に関するもので、AのB手段を具備する。…

出現回数	単語
18	
43	
60	
74	
31	
20	
	この発明は、半導体に…

【図 8】

【図 9】

出現回数	単語	正規化値
69	の	100.00
60	は	88.00
50		96.00
43	発明	94.00
42	手段	92.00
…	…	…
31	半導体	50.00
…	…	…

半導体製造

【図 10】

出現回数	単語	正規化値
125192	の	100.00
114837	に	88.04
103004	手段	88.52
…	…	…
5932	半導体	30.84

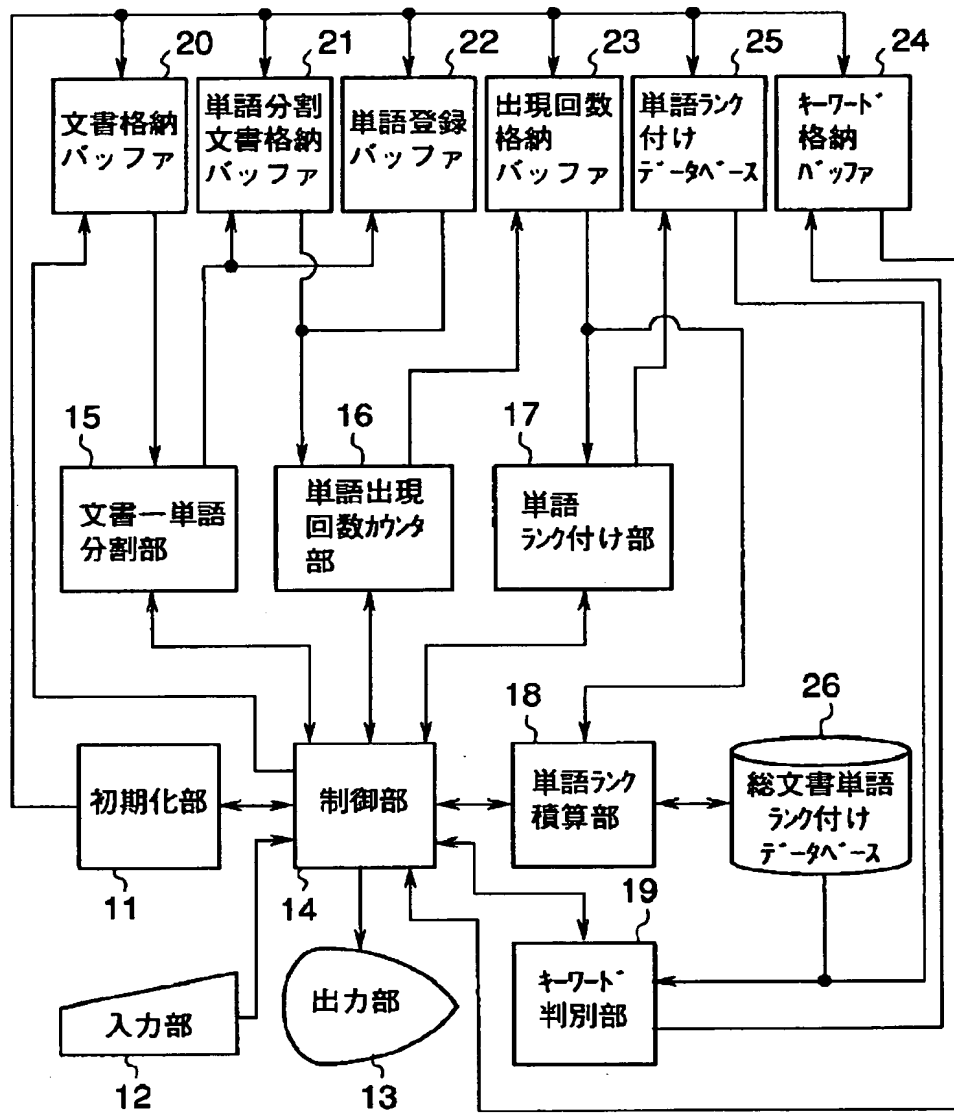
【図 13】

【図 17】

分野	キーワード
半導体	シリコン、ウェハ、IC、LSI、トランジスタ
メモリ	ROM、RAM、キャッシュ、PROM …
ネットワーク	LAN、WAN、WWW、インターネット…

文書番号
146
24
7
255
39

【図 2】



【図 18】

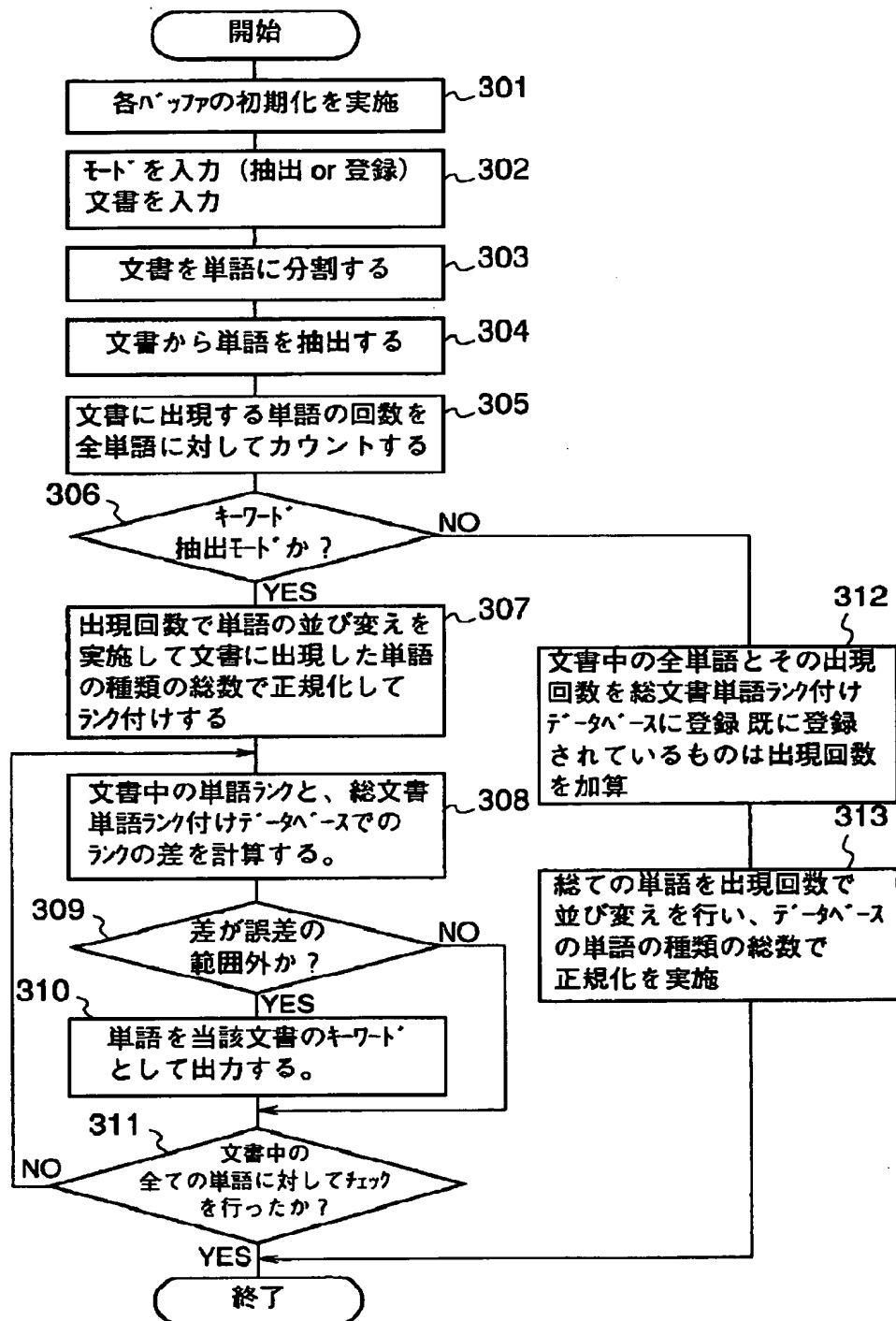
→ 文書番号

	1	2	3	4	...
この	1	1	1	1	...
発明	1	1	1	1	
、	1	1	1	1	
は	1	1	1	1	
半導体	1	0	0	1	
...					

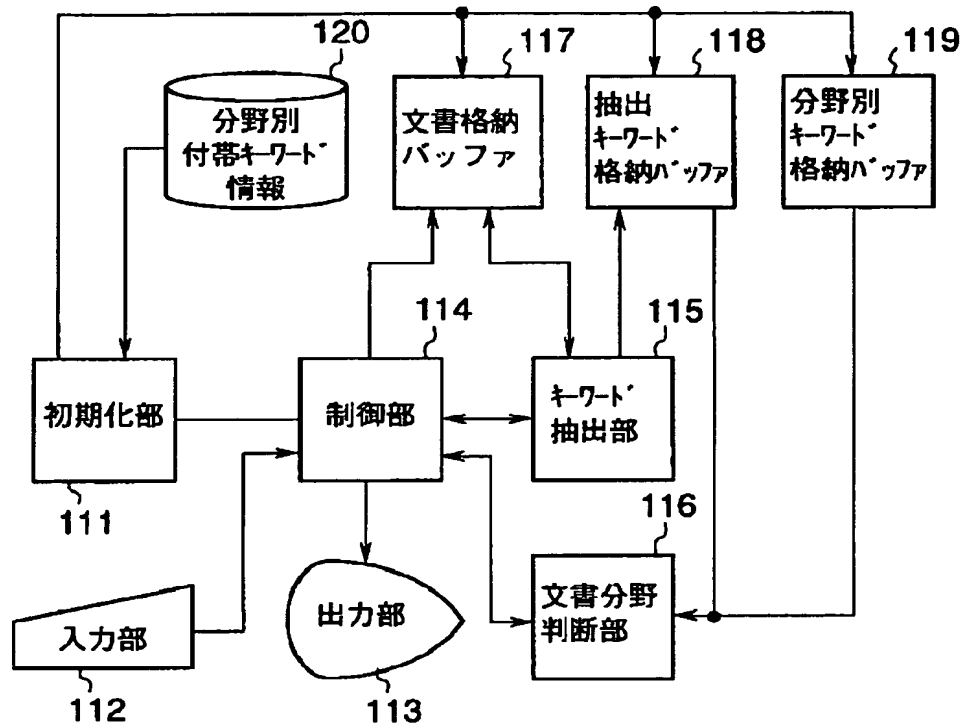
↑ 単語文字列

“1” は文書にその単語が存在していることを示す
 “0” は存在していないことを示す

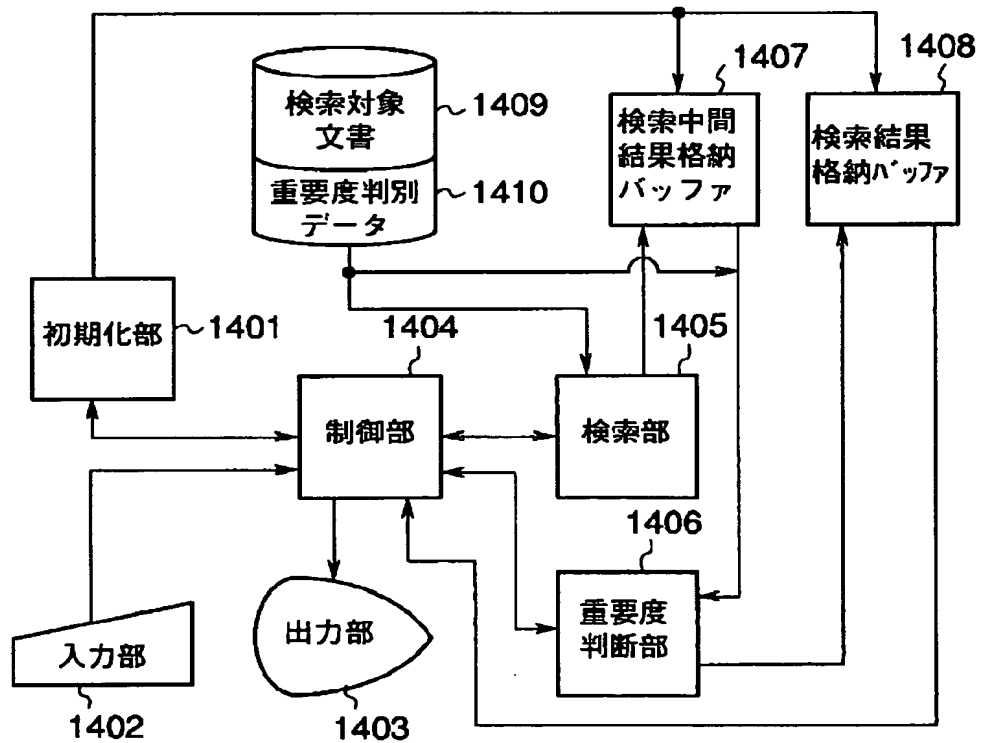
【図3】



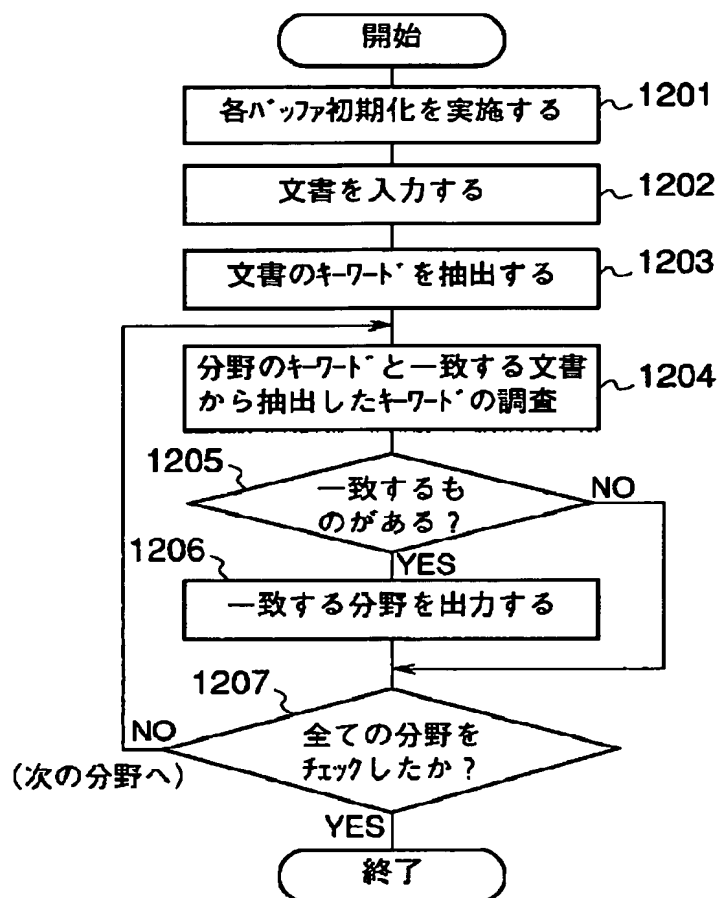
【図 11】



【図 14】



【図12】



【図15】

